

Guided macro-mutation in a graded energy based genetic algorithm for protein structure prediction

Mahmood A. Rashid^{a,d,*}, Sumaiya Iqbal^b, Firas Khatib^c, Md Tamjidul Hoque^b, Abdul Sattar^d

^a*School of Computing, Information and Mathematical Sciences, University of the South Pacific, Laucala Bay, Suva, Fiji*

^b*Department of Computer Science, University of New Orleans, LA, USA*

^c*Department of Computer and Information Science, University of Massachusetts Dartmouth, MA, USA*

^d*Institute for Integrated and Intelligent Systems, Griffith University, Brisbane, QLD, Australia*

Abstract

Protein structure prediction is considered as one of the most challenging and computationally intractable combinatorial problem. Thus, the efficient modeling of convoluted search space, the clever use of energy functions, and more importantly, the use of effective sampling algorithms become crucial to address this problem. For protein structure modeling, an off-lattice model provides limited scopes to exercise and evaluate the algorithmic developments due to its astronomically large set of data-points. In contrast, an on-lattice model widens the scopes and permits studying the relatively larger proteins because of its finite set of data-points. In this work, we took the full advantage of an on-lattice model by using a face-centered-cube lattice that has the highest packing density with the maximum degree of freedom. We proposed a graded energy—strategically mixes the Miyazawa-Jernigan (MJ) energy with the hydrophobic-polar (HP) energy—based genetic algorithm (GA) for conformational search. In our application, we introduced a 2×2 HP energy guided macro-mutation operator within the GA to explore the best possible local changes exhaustively. Conversely, the 20×20 MJ energy model—the ultimate objective function of our GA that needs to be minimized—considers the impacts amongst the 20 different amino acids and allow searching the globally acceptable conformations. On a set of benchmark proteins, our proposed approach outperformed state-of-the-art approaches in terms of the free energy levels and the root-mean-square deviations.

Keywords: *Ab initio* protein structure prediction; Genetic algorithms; FCC lattice; Miyazawa-Jernigan model; Hydrophobic-Polar model

*Corresponding author

Email addresses: mahmood.rashid@usp.ac.fj (Mahmood A. Rashid), siqbal1@uno.edu (Sumaiya Iqbal), fkhatib@umassd.edu (Firas Khatib), thoque@uno.edu (Md Tamjidul Hoque), a.sattar@griffith.edu.au (Abdul Sattar)

1. Introduction

Protein folding, by which the primary protein chain with amino acid residue sequence folds into its characteristics and functional three-dimensional (3D) structure in nature, is yet a very complex physical process to simulate [1, 2, 3]. Once the folded 3D shape is available, it enables protein to perform specific tasks for living organisms. Conversely, misfolded proteins are responsible for various fatal diseases, such as prion disease, Alzheimers disease, Huntingtons disease, Parkinsons disease, diabetes, and cancer [4, 5]. Because of these, protein structure prediction (PSP) problem has emerged as a very important research problem.

Homology modeling, threading and *ab initio* are the broad categories of available computational approaches. However, while homologous template is not available, *ab initio* becomes the only computation approach, which aims to find the three dimensional structure of a protein from its primary amino acid sequence alone such that the total interaction energy among the amino acids is minimized.

Ab initio computational approach for PSP is a daunting task [6] and for modeling the structure on a realistic continuum space such as off-lattice space is even more daunting. However, there are several existing off-lattice models such as Rosetta [7], Quark [8], I-TASSER [9], and so on which map the structures on the realistic continuum spaces rather than using discretized on-lattice spaces and hence, those approaches need to deal with the astronomical data-points incurring heavy computational cost. On-lattice model on the other hand, *i)* due to reduced complexity helps fast algorithms developments and *ii)* widens the scope as well as permits relatively longer protein chains to examine, which is otherwise prohibitive [10, 11, 12, 13, 14, 15]. The computed on-lattice fold can be translated to off-lattice space via hierarchical approaches to provide output in real-space [16, 17, 18, 19]. The Monte Carlo (MC) or, Conformational Space Annealing (CSA) used in Rosetta can be replaced with better algorithm developed using on-lattice models [16, 17, 18]. For instance, we embedded one of our previous on-lattice algorithm [17] within Rosetta and the embedded algorithm improved [20] the average RMSD by 9.5% and average TM-Score by 17.36% over the core Rosetta [7]. Similarly, the embedded algorithm also outperformed [20] I-TASSER [9]. These improvements motivated us further developing superior algorithms using on-lattice models.

The two most important building blocks of an *ab initio* PSP are *i)* an accurate (computable) energy function [19] and *ii)* an effective search or sampling algorithm. For a simplified model based PSP, it is possible to compute the lower bound [21]. It is also possible to know what would be the best score and hence the native score of a sequence by exhaustive enumeration [22, 23] (which is feasible to compute for smaller sequences only). Even though, there exists no efficient sampling algorithm yet that can conveniently obtain the known final structure starting from a random structure for all possible available cases [24, 25]. Therefore, a number of efforts are being made, such as, different types of meta-heuristics have been used in solving the on lattice PSP problems. These include Monte Carlo Simulation [26], Simulated Annealing [27], Genetic Algorithms (GA) [24, 25, 28, 29], Tabu Search with GA [30], GA with twin-removal operator [31], Tabu Search with Hill Climbing [32], Ant Colony Optimization [33], Particle Swarm Optimization [34, 35],

Immune Algorithms [36], Tabu-based Stochastic Local Search [37, 38], Firefly Algorithm [39], and Constraint Programming [40, 41].

Krasnogor et al. [42] applied HP model for PSP problem using the square, triangular, and diamond lattices and further extended their work applying fuzzy-logic [43]. Islam et al. further improved the performance of memetic algorithms in a series of work [44, 45, 46, 47] for the simplified PSP models. They also proposed a clustered architecture for the memetic algorithm with a scalable niching technique [48, 49, 50] for PSP. However, using 3D FCC lattice points, the recent state-of-the-art results for the HP energy model have been achieved by genetic algorithms [51, 52], local search approaches [38, 53], a local search embedded GA [54], and a multi-point parallel local search approach [55]. Kern and Lio [56] applied hydrophobic-core guided genetic operator for efficient searching on HP, HPNX and hHPNX lattice models. Several approaches towards the 20×20 energy model include a constraint programming technique used in [57, 58] by to predict tertiary structures of real proteins using secondary structure information, a fragment assembly method [59] to optimize protein structures. Among other successful approaches, a population based local search [60] and a population based genetic algorithm [61] are found in the literature that applied empirical energy functions.

In a hybrid approach, Ullah et al. [62] applied a constraint programming based large neighborhood search technique on top of the output of COLA [63] solver. The hybrid approach produced the state-of-the-art results for several small sized (less than 75 amino acids) benchmark proteins. In another work, Ullah et al. [64] proposed a two stage optimization approach combining constraint programming and local search using Berrera *et al.* [11] deduced 20×20 energy matrix (we denote this model as BM). In a recent work [65], Shatabda et al. presented a mixed heuristic local search algorithm for PSP and produced the state-of-the-art results using the BM model on 3D FCC lattice. Although the heuristics themselves are weaker than the BM energy model, their collective use in the random mixing fashion produce results better than the BM energy itself. In a previous work [66], we applied BM and HP energy models in a mixed manner within a GA framework and showed that hybridizing energies performs better than their individual performances.

In this work, we propose a graded as well as hybrid energy function with a genetic algorithm (GA) based sampling to develop an effective *ab initio* PSP tool. The graded energy-model strategically mixes 20×20 Miyazawa-Jernigan (MJ) contact-energy [10, 11] with the simple 2×2 Hydrophobic-Polar (HP) contact-energy model [12], denoted as MH (MJ+HP \rightarrow MH) in this paper. Specifically, we propose a hydrophobic-polar categorization of the HP model within a hydrophobic-core directed macro-mutation operator to explore the local benefits exhaustively while the GA sampling is guided by the MJ energy Matrix globally. While the fine grained details of the high resolution interaction energy matrix can become computationally prohibit, a low resolution energy model may effectively sample the search-space towards certain promising directions particularly emphasizing on the pair-wise contributions with large magnitudes which we have implementation strategically via a macro mutation. Further, we use an enhanced genetic algorithm (GA) framework [51] for protein structure optimization on 3D face-centered-cube (FCC) lattice model. Prediction in the FCC lattice model can yield the densest protein core [24] and the FCC lattice model can provide the maximum degree of freedom as well as the closest resemblance to the real or, high resolution folding within the lattice

constraint. FCC orientation can therefore align a real protein into the closest conformation amongst the available lattice configurations [25].

On a set of standard benchmark proteins, our MH model guided GA, named as *MH_GeneticAlgorithm* (MH_GA), shows significant improvements in terms of interaction energies and root-mean-square deviations in comparison to the state-of-the-art search approaches [62, 65, 61] for the lattice based PSP models. For a fair comparison, we run [62] and [65] using MJ energy model and in the result section, we compare our experimental results with the results produced by [62] and [65]. Further, we present an experimental analysis showing the effectiveness of using the hydrophobic polar categorization of the HP model to direct macro-mutation operation.

2. Background

Anfinsen’s hypothesis [67] and Levinthal’s paradox [68] form the basis and the confidence of the *ab initio* approach, which inform that the protein structure prediction can be relied only on the amino acid sequence of the target protein as well as there should be a non-exhaustive pathway to obtain the native fold. Thus, we set our goal to model the folding process using on-lattice model. Further, it has been argued in [69, 70], “... protein folding mechanisms and landscapes are largely determined by the topology of the native state and are relatively insensitive to details of the interatomic interactions. This dependence on low resolution structural features, rather than on high resolution detail, suggests that it should be possible to describe the fundamental physics of the folding process using relatively low resolution models ... The observation that protein folding mechanisms are determined primarily by low resolution topological features and not by high resolution details suggests that a simple theory incorporating features of the native state topology should be successful in predicting the broad outlines of folding reactions”. A rigorous discussion on on-lattice models can be found in [71, 72]. Further, exploration of an astronomically large search space and the evaluation of the conformations using a real energy models are the big challenges often being computationally prohibitive, especially for sequences > 200 residues long, whereas simplified models can aid in modeling and understanding the protein folding process feasibly. Next, we describe the model that we use.

2.1. Simplified model

In our simplified model, we use 3D FCC lattice points to map the amino acids of a protein sequence. In the mapping, each amino acid of the sequence, occupies a point on the lattice to form a continuous chain of a self-avoiding-walk. We apply the MJ energy matrix in conjunction with the HP energy model in a genetic algorithm framework for PSP. The FCC lattice, the HP and MJ energy models, and the GA are briefly described below.

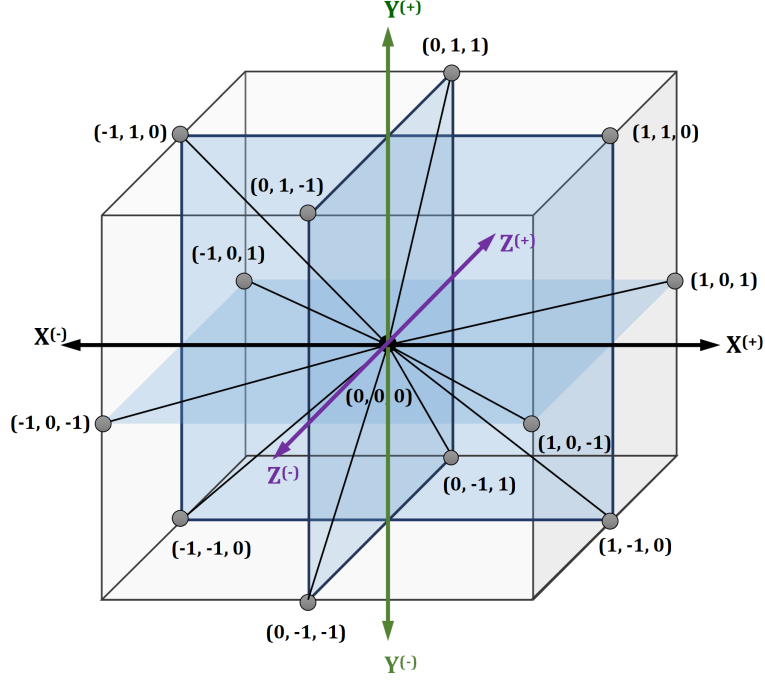


Figure 1: A 3-dimensional face-centered-cubic lattice space. The 12 basis vectors of the neighbors of the origin $(0,0,0)$ in a Cartesian coordinate system.

FCC lattice

The FCC lattice has the highest packing density compared to the other existing lattices [73]. Thus, FCC model can provide maximum degree of freedom within a constrained space. In FCC, each lattice point (the origin in Figure 1) has 12 neighbors with closest possible distance having 12 *basis vectors* as follows:

$$\begin{aligned}
 v_1 &= (1, 1, 0) & v_4 &= (-1, -1, 0) & v_7 &= (-1, 1, 0) & v_{10} &= (0, 1, -1) \\
 v_2 &= (1, 0, 1) & v_5 &= (-1, 0, -1) & v_8 &= (1, -1, 0) & v_{11} &= (1, 0, -1) \\
 v_3 &= (0, 1, 1) & v_6 &= (0, -1, -1) & v_9 &= (-1, 0, 1) & v_{12} &= (0, -1, 1)
 \end{aligned}$$

In simplified PSP, conformations are mapped on the lattice by a sequence of basis vectors, or by the *relative vectors* that are relative to the previous basis vectors in the sequence.

HP energy model

Based on the hydrophobic property, the 20 amino acids which are the constituents of all proteins, are broadly divided into two categories: (a) hydrophobic amino acids (*Gly, Ala, Pro, Val, Leu, Ile, Met, Phe,*

Tyr, *Trp*) are denoted as H; and (b) hydrophilic or polar amino acids (*Ser*, *Thr*, *Cys*, *Asn*, *Gln*, *Lys*, *His*, *Arg*, *Asp*, *Glu*) are denoted as P. In the 2×2 HP model [12], when two non-consecutive hydrophobic amino acids become topologically neighbors, they contribute a certain amount of negative energy, which for simplicity is considered as -1 (Table 1). The total energy E_{HP} (Equation 1) of a conformation based on the HP model becomes the sum of the contributions over all pairs of the non-consecutive hydrophobic amino acids (Figure 2a).

Table 1: The 2×2 HP energy model.

	H	P
H	-1	0
P	0	0

$$E_{HP} = \sum_{i < j-1} c_{ij} \times e_{ij} \quad (1)$$

where, $c_{ij} = 1$ if amino acids at positions i and j in the sequence are non-consecutive but topological neighbors on the lattice, otherwise $c_{ij} = 0$. The $e_{ij} = -1$ if the i th and j th amino acids are both hydrophobic, otherwise $e_{ij} = 0$.

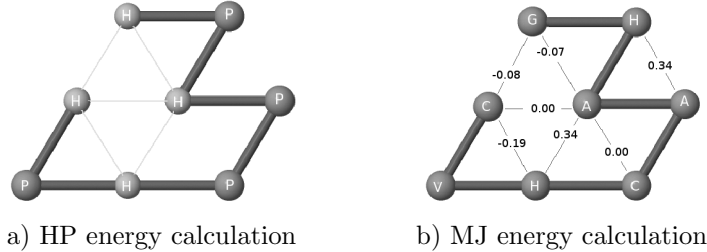


Figure 2: On a lattice based model, (a) is showing H-H contact energy -4 (-1×4) using the HP model for a random sequence: *HPHPHPHP* (b) is showing the sum of the pair-wise contact potentials 0.34 ($(-0.08) + (-0.19) + (0.00) \times 2 + (-0.07) + (0.34) \times 2$) using the MJ model for a random sequence: *GHAACHVC*.

MJ energy model

By analyzing crystallized protein structures, Miyazawa and Jernigan [10] statistically deduced a 20×20 energy matrix (better known as MJ energy model) that considers residue contact propensities between the amino acids. BM is a similar energy matrix as MJ deduced by Berrera *et al.* [11] by calculating empirical contact energies on the basis of information available from a set of selected protein structures and following the quasi-chemical approximation. In this work, we use MJ energy model. The total energy E_{MJ} (Equation

2.2. Genetic algorithms

GAs [74] are a family of population-based search algorithms which can be applied for PSP as an optimization problem. The outline of GA as given in Algorithm 1, follows simple steps: Line 1 initializes the population; the Line 2 evaluates the solutions to rank them by relative quality; and the Lines 4-7 are repeating on generating, evaluating and replacing the least-fitted off-springs within the population until the termination criteria arises. For the coding scheme, non-isomorphic encoding [75] has been applied and the v_1, \dots, v_{12} (in Figure 1) can be thought of be renamed as a, \dots, l respectively.

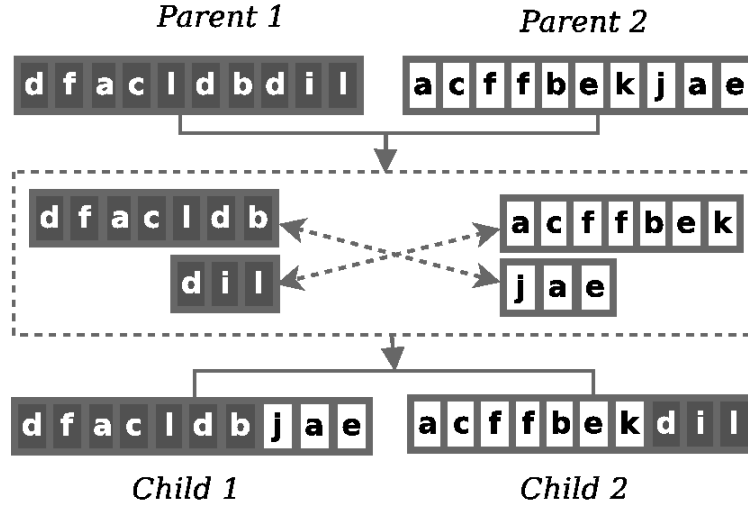


Figure 3: Typical crossover operator: exchanging parts and forming new chromosomes.

A typical crossover operator randomly splits two solutions at a randomly selected crossover point and exchanges the parts between them (Figure 3) and a typical mutation operator alters a solution at a random point (Figure 4). In the case of PSP, conformations are regarded as solutions of a GA population.



Figure 4: Typical mutation operator: mutating one point into some other point.

3. Methods

This section describes the proposed MH_GA framework along with the implementation level detail. We implemented the framework in Java (J2EE), using Rocks clusters. The code for MH based GA is freely

available online¹.

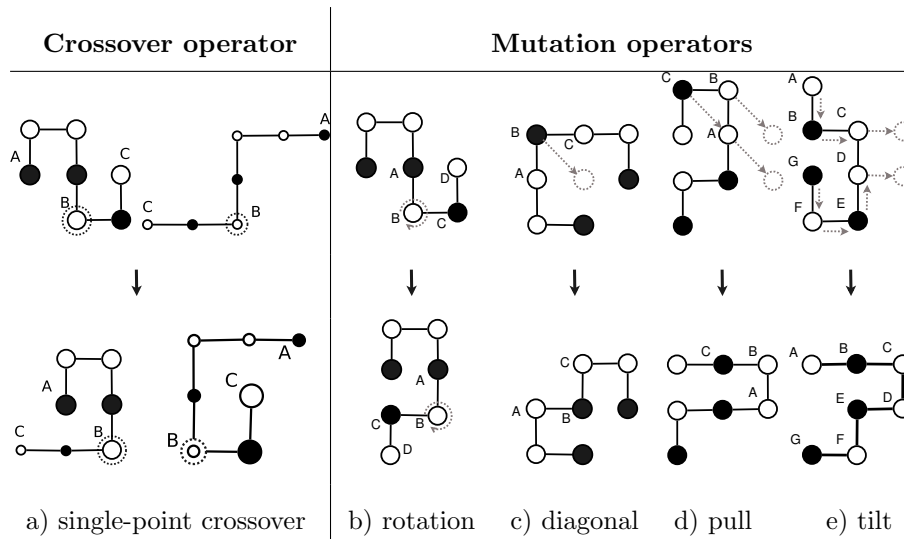


Figure 5: The primitive operators used in our genetic algorithms. The crossover operator applied on two parent conformations to exchange their parts to generate two child conformations (as shown in a) and the mutation operators are applied on single conformation to generate single child conformation (as shown in b, c, d, and e). The operators are implemented in 3D space, however, for simplification and easy understanding the figures are presented in 2D space. The black solid circles represent the hydrophobic amino acids and others are polar.

3.1. The primitive operators implemented in the GA framework

The primitive operators that we implemented within the MH_GA framework are crossover (Figure 5a), rotation mutation (Figure 5b), diagonal move (Figure 5c), pull moves (Figure 5d), and tilt moves (Figure 5e). The Rotation, diagonal move, pull moves and tilt moves are implemented as mutation operators.

1. **Crossover:** At a given crossover point (dotted circle in Figure 5a), two parent conformations exchange their parts and generate two children. The success rate of crossover operator decreases with the increase of the compactness of the structure.
2. **Rotation:** One part of a given conformation is rotated around a selected point (Figure 5b). This move is mostly effective at the beginning of the search.
3. **Diagonal move:** Given three consecutive amino acids at lattice points A , B , and C , a diagonal move at position B takes the corresponding amino acid diagonally to a free position (Figure 5c). The diagonal moves are very effective on FCC lattice [37, 41] points.

¹Download JAR files from: http://cs.uno.edu/~tamjid/Software/MH_GA/JarFiles.zip

4. **Pull moves:** The amino acids at points A and B are pulled to the free points (Figure 5d) and the connected amino acids are pulled as well to get a valid conformation. The pull moves [23] are local, complete, and reversible. These are very effective especially when the conformation is compact.
5. **Tilt moves:** Two or more consecutive amino acids connected in a straight line are moved by a tilt move to immediately parallel lattice positions [29]. The tilt-moves pull the conformation from both sides until a valid conformation is found. In Figure 5e, the amino acids at points C and D are moved and subsequently other amino acids from both sides are moved as well.

3.2. Genetic algorithm framework

The *pseudocode* of MH_GA framework is presented in Algorithm 2. It uses a set of primitive operators (Figure 5) in an exhaustive generation approach to diversify the search, a hydrophobic core-directed macro-mutation operator to intensify the search, and a random-walk algorithm to recover from the stagnation. Like other search algorithms, it requires initializing the population and the solutions need to be evaluated in each iteration.

Algorithm 2: MH_GeneticAlgorithm (MH_GA)

```

/* INPUT: Protein/Amino acid sequence, popSize: Population size; opR: Operator selection probabilities
*/
/* OUTPUT: Global best conformation */
/* VARIABLE: op: Operators; c, c': Conformations; curP, newP: Current and new populations; mmCount:
Macro-mutation counter; rwCount: Non-improving random-walk counter */

1 curP ← initialise ;
2 repeat
3   op ← selectOperator (opR);
4   if (op is crossover) then
5     /* ** go for crossover */
6     while (¬ full (newP)) do
7       c, c' ← randomConfs (curP);
8       newP.add (doCrossover);
9     end
10  else if (op is mutation) then
11    /* ** go for mutation */
12    foreach (c ∈ curP) do
13      newP.add (doMutation);
14    end
15  else
16    /* ** go for macro-mutation */
17    foreach (c ∈ curP) do
18      newP.add (doHCDMacroMutation);
19    end
20  end
21  if (¬ improved (newP, rwCount)) then
22    newP ← goRandomWalk ;
23  end
24  curP ← newP;
25 until (termination criteria);
26 return bestConformation (curP);

```

The algorithm initializes (Algorithm 2: Line 7) the current population with randomly generated individuals. At each generation, it selects a genetic operator based on a given probability distribution to use through the generation (Algorithm 2: Line 9). In fact, we select the operators randomly by giving equal opportunities to all operators. The selected operator is used in an exhaustive manner (Algorithm 2: Line 11-12 or Line 14-16) to obtain all conformations in the new population. We ensure that no duplicate conformation is added to the new population. The `add()` method (Line 12 or 16 in Algorithm 2) takes care of adding the non-duplicate conformations to the new population. For a given number of generations, if the best conformation in the new population is not better than the best in the current population, our algorithm triggers a random-walk technique (Algorithm 2: Line 18) to diversify the new population. Nevertheless, after each generation, the new population becomes the current population (Algorithm 2: Line 19); and the search continues. Finally, the best conformation found so far is returned (Algorithm 2: Line 20). Along with MJ potential matrix, the HP energy model is used during move selection by the macro-mutation operator. The macro-mutation operator is used as other mutation operators (Figure 5b-e) in MH-GA. The details of initialization, evaluation of fitness, exhaustive generation, macro-mutation and stagnation recovery schemes are presented below.

Initialization

Our algorithm starts with a feasible set of conformation known as population. We generate initial conformations following a self-avoiding walk on FCC lattice points. The *pseudocode* of the algorithm is presented in Algorithm 3. It places the first amino acid at $(0, 0, 0)$. It then randomly selects a basis vector to place the successive amino acid at a neighboring free lattice point. The mapping proceeds until a self-avoiding walk is found for the whole protein sequence.

Evaluate the fitness

For each iteration, the conformation is evaluated by calculating the contacts (topological neighbor) potentials where the two amino acids are non-consecutive. The pseudo-code in Algorithm 4 presents the algorithm for calculating the interaction energy of a given conformation. The contact potentials are found in MJ potential matrix [10] (*see* Table 2).

Exhaustive generation

Unlike standard genetic algorithm, in MH-GA, the randomness is reduced significantly by applying exhaustive generation approach. For mutation operators, MH-GA adds one resultant conformation to the new population that corresponds to *each* conformation in the current population. Operators are applied to all possible point (Algorithm 5) exhaustively until finding a better solution than the parent. If no better solution is found, the parent survives through the next generation. On the other hand, for crossover operators, two resultant conformations are added to the new population from two randomly selected parent conformations. Crossover operators generate child conformations by applying the crossover operator in all possible points

Algorithm 3: initialise

```
/* Is called from Algorithm 2 in Line 1 */
/* INPUT: Protein/Amino acid sequence, FCC basis vectors, popSize: Population size */
/* OUTPUT: Initial population */
/* VARIABLE: AA: Array of amino acid; c: Conformations; point: Unoccupied point on 3D FCC Lattice space */
1 for (p = 1; p ≤ popSize; p++) do
2   AA[0] ← aminoAcid (0,0,0);
3   for a number of times do
4     for (i = 1 to seqLength-1) do
5       j ← getRandom (12);
6       point ← AA[i - 1] + basisVector[j];
7       if point is not free then
8         break;
9       else
10        AA[i] ← aminoAcid (point);
11      end
12    end
13  end
14  if full structure found then
15    c.AminoAcid ← AA [];
16  else
17    c ← a deterministic structure;
18  end
19  c.fitness ← evaluate (c.aminoAcid);
20  initPop.add (c);
21 end
22 return initPop
```

Algorithm 4: evaluate

```
/* Is called from Algorithm 3 in Line 19 */
/* INPUT: MJ energy matrix(20 × 20), AA: Array of amino acid */
/* OUTPUT: Fitness of the structure */
/* VARIABLE: seqLength: Sequence length; pointl, pointj: Occupied point on 3D FCC Lattice space */
1 fitness ← 0
2 for (i = 0 to seqLength - 1) do
3   for (j = i + 2 to seqLength - 1) do
4     pointl ← AA [i];
5     pointj ← AA [j];
6     sqrD ← getSqrDist (pointl, pointj);
7     if sqrD = 2 then
8       fitness ← fitness + Ebm[i][j];
9     end
10  end
11 end
12 return fitness;
```

(Algorithm 6) on two randomly selected parents. The best two conformations from the parents and the children are then become the resultant conformations for the next generation.

Algorithm 5: doMutation

```
/* Is called from Algorithm 2 in Line 11 */  
/* INPUT: conf : Conformation */  
/* OUTPUT: Best mutated conformation */  
/* VARIABLE: c : Conformation; offspring : List of type conformation */  
1 offspring.add(conf);  
2 foreach ( $1 \leq \text{pos} \leq \text{seqLength}$ ) do  
3    $c \leftarrow \text{applyOperator}(\text{conf}, \text{pos})$ ;  
4   offspring.add(c);  
5 end  
6 return bestConformation(offspring);
```

Algorithm 6: doCrossover

```
/* Is called from Algorithm 2 in Line 7 */  
/* INPUT:  $c_1$  and  $c_2$  : Conformations */  
/* OUTPUT: Best two conformations after crossover */  
/* VARIABLE:  $c, c'$  : Conformations; offspring : List of type conformation */  
1 offspring.add( $c_1, c_2$ );  
2 foreach ( $1 \leq \text{pos} \leq \text{seqLength}$ ) do  
3    $c, c' \leftarrow \text{applyOperator}(c_1, c_2, \text{pos})$ ;  
4   offspring.add( $c, c'$ );  
5 end  
6 return best2Conformations(offspring);
```

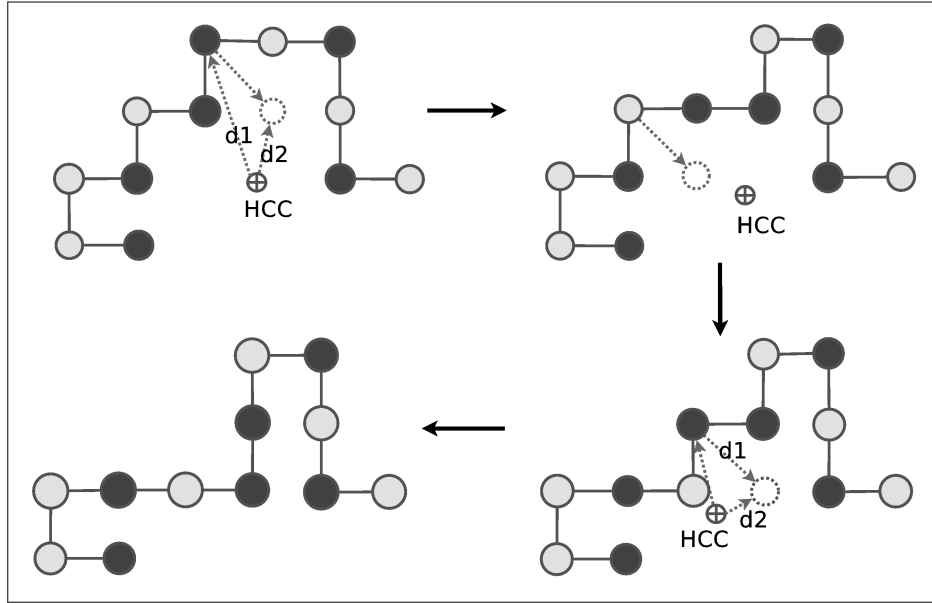


Figure 6: A macro-mutation operator repeatedly used diagonal moves. The moves of an amino acid are guided by the distance of current position (d_1) and the distance of target position (d_2) from the HCC. The operator is implemented in 3D space, however, for simplification and easy understanding, the figures are drawn in 2D space.

Macro-mutation operator

Protein structures have hydrophobic cores (H-Core) that hide the hydrophobic amino acids from water and expose the polar amino acids to the surface to be in contact with the surrounding water molecules [76]. H-core formation is an important objective of HP based PSP. Macro-mutation operator is a composite operator (Figure 6) that uses a series of diagonal-moves (Figure 5c) on a given conformation to build the H-core around the hydrophobic-core-center (HCC). The macro-mutation squeezes the conformation and quickly forms the H-core. In MH-GA, macro-mutation is used as other mutation operators. The Algorithm 7 presents the *pseudocode* of macro-mutation algorithm.

Algorithm 7: doHCDMacroMutation

```

/* Is called from Algorithm 2 in Line 15                                     */
/* INPUT: HP energy matrix( $2 \times 2$ ), C: Conformation; repeat: Loop counter      */
/* OUTPUT: Mutated conformation                                           */
/* VARIABLE: T: Either hydrophobic (H) or polar (P); AA: Array of amino acid */
1 AA [ ]  $\leftarrow$  C.AminoAcid[ ];
2 for i = 1 to repeat do
3   T  $\leftarrow$  P if bernoulli( $p$ ), else H
4   AA[j] : jth amino acid in conformation
5   point: unoccupied new position for AA[j]
6   hcc  $\leftarrow$  findHCC ()
7   foreach j : typeOf(AA[j]) = T do
8     dold  $\leftarrow$  getDistance (AA[j],hcc)
9     if T = P then
10      point  $\leftarrow$  findFreePoint (AA[j])
11      applyDiagonalMove(AA[j], point)
12    end
13    else
14      point  $\leftarrow$  findFreePoint (AA[j])
15      dnew  $\leftarrow$  getDistance (point,hcc)
16      if dnew  $\leq$  dold then
17        applyDiagonalMove(AA[j], point)
18        break
19      end
20    end
21  end
22 end
23 C.AminoAcid[ ]  $\leftarrow$  AA [ ];
24 return C

```

In macro-mutation, the HCC is calculated by finding arithmetic means of x , y , and z coordinates of all H amino acids. In macro-mutation, for a given number of iterations, diagonal moves apply repeatedly either at each P- or at each H-type amino acid positions. Whether to apply the diagonal move on P- or H-type amino acids is determined by using a *Bernoulli* distribution (Algorithm 7 : Line 2) with probability p (intuitively we use $p = 20\%$ for P-type amino acids). For a P-type amino acid, the first successful diagonal move is considered. However, for a H-type amino acid, the first successful diagonal move that does not increase the Cartesian distance of the amino acid from the HCC is taken. All the amino acids are traversed and the successful moves are applied as one composite move.

3.3. Stagnation recovery

Like other search algorithm, GA can get stuck in the local minima or, can be stalled. Stall condition can occur when similarities with the chromosomes in GA increases heavily and the operators are unable to produce better diverse solutions. Further, with the PSP search, resulting solutions become phenotypically compact which reduce the likelihood of producing better solution from the population due to harder self-avoid-walk (SAW) constraints [17, 77, 78]. It would rather require very intelligent moves to reform into another competitive compact SAW. To deal with such situation, we apply the following two actions:

Removing duplicates

In genetic algorithm it has been observed that with increasing generations, the similarity among the individuals within the population increases. In worst case scenario, all the individuals become similar and forces the search to stall in the local minima. In our approach, we remove duplicates from each generation to maintain the diversity of the population. During exhaustive generation, we check the existence of the newly generated child in the new population. If it does not exist then the new solution is added to the new population list. Our approach reduces the frequency of stagnations.

Applying random-walk

Sometimes, early convergence leads the search towards the stagnation situation. In the HP energy model, premature H-cores are observed at local minima. To break these H-cores, in MH-GA (Algorithm 2 : Line 18), a random-walk algorithm (Algorithm 8) is applied. This algorithm uses pull moves [23] (as shown in Figure 5d) to break the H-core. We use pull-moves because they are complete, local, and reversible. Successful pull moves never generate infeasible conformations. During pulling, energy level and structural diversification are observed to maintain balance among these two. We allow energy level to change within 5% to 10% that changes the structure from 10% to 75% of the original. We try to accept the conformation that is close to the current conformation in terms of the energy level but as far as possible in structural diversity, and which is determined by the function `checkDiversity()` in Algorithm 8 at Line 5. For genetic algorithm, random-walk is very effective [79] to recover from stagnation.

The complete flow of MH-GeneticAlgorithm (Algorithm: 2) is graphically presented in Figure 7. Further, it describes the steps taken within macro mutation procedure (Algorithm: 7).

4. Performance Evaluation

To compare and evaluate the performance of the proposed PSP predictor with respect to the state-of-the-art approaches, we used the measures *Relative Improvement (RI)* and *RMSD* comparisons. They are defined below:

Algorithm 8: goRandomWalk

```

/* Is called from Algorithm 2 in Line 19                                     */
/* INPUT: inPop: Current population; pct: Changed percentage (%)          */
/* OUTPUT: New diverged population                                         */
/* VARIABLE: outPop: New diverged population;    AA: Array of amino acid; c,c': Conformations */
1 foreach (c ∈ inPop) do
2   isFound ← false;
3   AA [ ] ← c.AminoAcid [ ];
4   while (¬isFound) do
5     for (i = 1; i ≤ seqLength & ¬isFound; i++) do
6       applyPullMove(AA[i]);
7       c'.AminoAcid [ ] ← AA [ ];
8       isFound ← checkDiversity (c,c',pct);
9     end
10  end
11  outPop.add (c');
12 end
13 return outPop

```

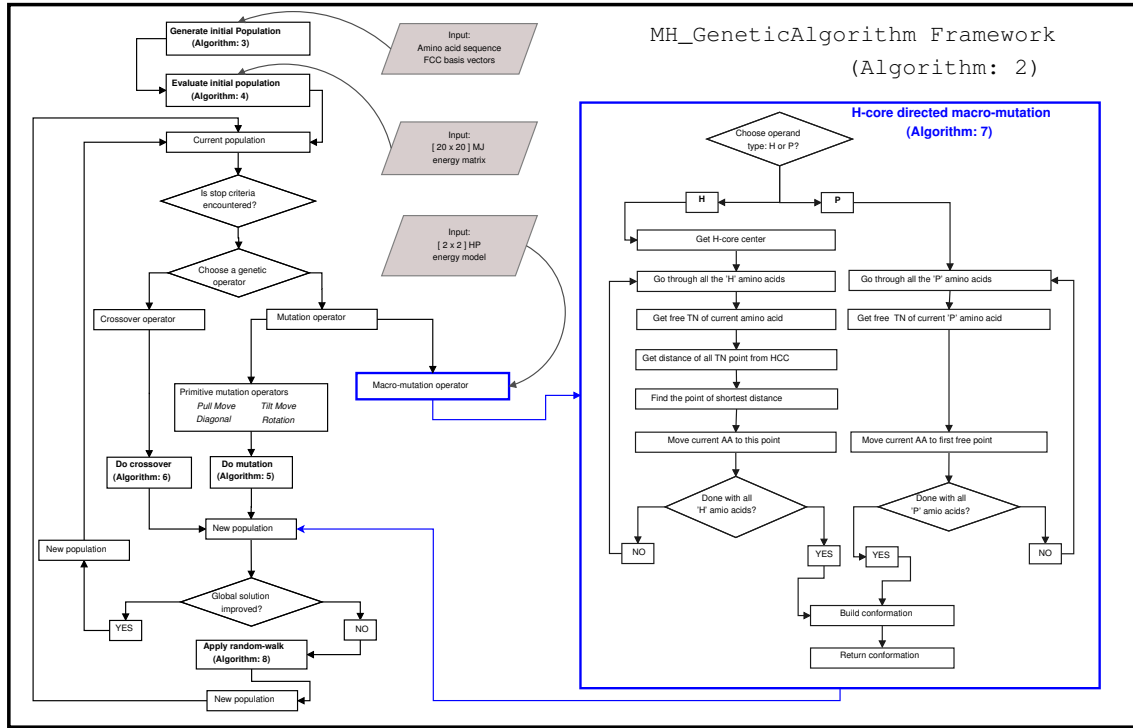


Figure 7: A complete overview of our algorithmic approach. The macro-mutation procedure is described step by step (inside the blue box). The procedural sub blocks are marked in bold along with the corresponding labels of the algorithms described above.

Relative Improvement (RI)

The difficulty to improve energy level is increased as the predicted energy level approaches to a known lower bound of a given protein. For example, if the lower bound of free energy of a protein is -100 , the efforts to improve energy level from -80 to -85 is much less than that to improve energy level from -95 to -100 though the change in energy is the same (-5). The RI computes the relative improvements that our algorithm (target, t) achieved w.r.t. the state-of-the-art approaches (reference, r).

For each protein, the relative improvement of the target (t) w.r.t. the reference (r) is calculated using the formula in Equation 3, where E_t and E_r denote the average energy values achieved by target and reference respectively.

$$\text{RI} = \frac{E_t - E_r}{E_r} \times 100\% \quad (3)$$

RMSD comparison

The root mean square deviation (RMSD) is frequently used to measure the differences between values predicted by a model and the values actually observed. We compare the predicted structures obtained by our approach with the state-of-the-art approaches by measuring the root-mean-square w.r.t. the native structures from PDB. For any given structure the root-mean-square is calculated using Equation 4,

$$\text{RMSD} = \sqrt{\frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n (d_{ij}^p - d_{ij}^n)^2}{n \times (n-1)/2}} \quad (4)$$

where d_{ij}^p and d_{ij}^n denote the distances between i th and j th amino acids respectively in the predicted structure and the native structure of the protein. The average distance between two α -Carbons in native structure is 3.8\AA . To calculate root-mean-square, the distance between two neighbor lattice points is considered as 3.8\AA .

5. Results and discussion

In this section, we discuss the obtained results along with the comparison of the performance of MH.GeneticAlgorithm with the other state-of-the-art results [61, 62, 65]. Further, we present an analysis of the results.

5.1. Benchmark

In our experiment, the protein instances are taken from the literatures. The first seven proteins (*4RXN*, *1ENH*, *4PTI*, *2IGD*, *1YPA*, *1R69*, and *1CTF*) in Table 3 are taken from [62] and [65], and the next five proteins (*3MX7*, *3NBM*, *CMQO*, *3MRO*, and *3PNX*) are taken from [65]. The two other protein instances in Table 5 (*2J61* and *2HFQ*) are taken from [61].

Table 3: The benchmark proteins used in our experiments.

ID	Len	Protein sequence
4RXN	54	MKKYTCTVCGYIYNPEDGDPDNGVNPGETDFKDIPDDWVCPLCGVGKDQFEEVEE
1ENH	54	RPRTAFSSEQLARLKREFNENRYLTERRRQQLSSELGLNEAQIKIWFQNKRAKI
4PTI	58	RPDFCLEPPYTGPCKARIIRYFYNAKAGLCQTFVYGGCAKRNNFKSAEDCMRTCGGA
2IGD	61	MTPAVTTYKLVLINGKTLKGETTTKAVDAETAEKAFKQYANDNGVDGVWYDDATKFTVTTE
1YPA	64	MKTEWPELVGKAVAAAKVILQDKPEAQIIVLPVGTIVTMEYRIDRVRLFVDKLDNIAQVPRVG
1R69	69	SISSRVKSKRIQLGLNQAELAQKVGTTQQSIEQLENGKTKRPRFLPELASALGVSVDWLLNGTSDSNVR
1CTF	74	AAEEKTEFDVILKAAGANKVAVIKAVRGATGLGLKEAKDLVESAPAALKEGVSKDDAEALKKALEEAGAEVEVK
3MX7	90	MTDLVAVWDVALSDGVHKEIEFEHGTTSKGRVVYVDGKEEIRKEWMFKLVGKETFYVGAAKTKATINIDAISGFA YEYTLINGKSLKKYM
3NBM	108	SNASKELKVLVLCAGSGTSAQLANAINEGANLTVRVIANSGAYGAHYDIMGVYDLIILAPQVRSYYREMKVDA ERLGIQIVATRGMEYIHLTKSPSKALQFVLEHYQ
3MQO	120	PAIDYKTAFLHAPIGLVLSRDRVIEDCNDELAEIFRCARADLIGRSFEVLYPSSDEFERIGERISPVMIAGSY ADDRIMKRAGGELFWCHVTGRALDRDTAPLAAGVWTFEDLSATRRVA
3MRO	142	SNALSASEERFQLAVSGASAGLWDWNPKTGAMVLSPHFKKIMGYEDHELPEITGHRESIHPDDRARVLAALKA HLEHRDITYVEYRVTRSGDFRWIQSRGQALWNSAGEPYRMVGWIMDVDRKREDAALRVSRRELRL
3PNX	160	GMENKKMNLFFSGDYDKALASLIIANAAREMEIEVTFICAFWGLLLLRDPEKASQEDKSLYEQAFFSLTPREA EELPLSKMNLGGIGKKMLLEMMKEEKAPKLSDLLSGARKKEVKFYACQLSVEIMGFKKEELFPEVQIMDVKEYL KNALESDQLFI
2J6A	135	MKFLTTNFKCSVKACDTSNDNFPLQYDGSKCQLVQDESIEFNPEFLNIVDRVDWPAVLTVAAELGNALPPT KPSFPSSIQELTDDMAILNDLHTLLQTSIAEGEMKCRNCGHYIYKNGIPNLLPPLHV
2HFQ	85	MQIHVYDITYKAKDGHVMHFDVFTDVRDDKKAIEFAKQWLLSSIGEAGATVTSEECRCHFSQKAPDEVIEAIKQN GYFIYKMEGCN
3MSE	180	GISPVLNMMKSYMKHSNIRNIIINIMAHELSVINNHKIYINELFYKLDTNHNGSLSHREIYTVLASVGKKWD INRILQALDINDRGNITYTEFMAGCYRWKNIESTFLKAAFNKIDKDEGDYISKSDIVSLVHDKVLDNNDIDNFF LSVHSIKKGIPREHIINKISFQEFQDYMSTF
3MR7	189	SNAERRLCAILAADMAGYSRLMERNETDVLNRQKLYRRELIDPAIAQAGGQIVKTTGDGMLARFDTAQAALRCA LEIQQAMQQREEDTPRKRIQYRIGINIGDIVLEDGDI FGDAVNVAARLEAISEPGAICVSDIVHQITQDRVSE PFTDLGLQKVKNITRPIRVWQWVPDADRDQSHDPQPSHVQH
3MQZ	215	SNAMSVQTIERLQDYLLPEWVSIFDIADFSGRMLRIRGDIRPALLRLASRLAELLNESGPRPWYPHVAHMRRR VNPPPETWLAGPEKRGYKSYAHSGVFIGRGLSVRFILKDEATEERKNLGRWMSRSGPAFEQWKKKVGDLRDFG PVHDDPMADPPKVEWDPRVFGERLGSLSASLDIGFRVTFDTSLAGIVKTI RTFDLLYAEAEKGS
3NO3	238	GKDNKTVIAHRYGKTEGSAQNSIRSLERASEIGAYGEFVHLTADNVLVYHDNDIQGKHIQSCYDELKDLQ LSNGEKLPTLEQYLKRAKKLNIRLIFELKSHDTPERNRDAARLSVQMVKRMKLAKRTDYISFNMDACEFIRLC PKSEVSYLNGELSPMELKELGFTGLDYHYKVLQSHPDWVKDCKVLGMTSNVWTVDDPKLMEEMIDMGVDFITDDL PEETQKILHSRAQ
3NO7	248	MGSDKIHHHHHENLYFQGMTFSKELREASRPIDDIYNDGFIQDLLAGKLSNQAVRQYLADASYLKEFTNIYA MLIPKMSSMEDVKFLVEQIEFMLEGEVEAEVLADF INEPYEEIVKEKVWPPSGDHYIKHMYNFAFARENAAFTI AAMAPCPYVYAVIGKRAMEDPKLNKESVTSKWFQFYSTEMDELVDVFDQLMDRLTKHCSETEKKEIKENFLQSTI HERHFFNMAYINEKWEYGGNNNE
3ON7	280	GMKLETIDYRAADSARFVESLRETGFGVLSNHPIDKELVERIYTEWQAFFNSEAKNEFMFNRETHDGFFPASIS ETAKGHTVKDIKEYHYVYPWGRIPDSLRANILAYYEKANTLASELLEWETYSPEIKAKFSIPLPEMIANSHKT LLRILHYPPMTGDEEMGAIRAAAHEDINLITVLPANEPGLQVKAKDGSWLDVPSDFGNIIINIGDMLQEASDGY FPSTSHRVINPEGTDKTKSRISLPLFLPHPSVVLSEYRTADSYLMERLRELGL

5.2. Comparing with the state-of-the-art

In the literature we found very few works [60, 61] that used 20×20 MJ potential-matrix [10] for protein structure prediction on 3D FCC lattice. However, Torres *et al.* [61] used 3D HCP lattice and Kapsokalivas *et al.* [60] used 3D cubic lattice in their works for protein mapping. In other works, Ullah *et al.* [62] and

Table 4: The energy values are obtained from different algorithms for the specified energy models. The average values are calculated over 50 different runs. The bold-faced values indicate the winner (the lower the better).

Protein details			The state-of-the-art						Our approach			
			Hybrid [62]			Local Search [65]			The MH_GA			
			MJ energy		Time	MJ energy		Time	MJ energy		Time	RI
<i>Seq</i>	<i>Size</i>	<i>H</i>	<i>Best</i>	<i>Avg</i>	<i>Avg</i>	<i>Best</i>	<i>Avg (r)</i>	<i>Avg</i>	<i>Best</i>	<i>Avg (t)</i>	<i>Avg</i>	over [65]
4RXN	54	27	-32.61	-30.94	1:02:12	-33.33	-31.21		-36.36	-33.60		7.66%
1ENH	54	19	-35.81	-35.07	1:02:03	-29.03	-28.18		-38.39	-35.67		26.58%
4PTI	58	32	-32.07	-29.37	1:01:26	-31.16	-28.33		-35.65	-31.01		9.46%
2IGD	61	25	-38.64	-32.54	1:43:08	-32.36	-28.29	1:00:00	-36.49	-33.75	1:00:00	19.30%
1YPA	64	38	<i>n/a</i>	<i>n/a</i>		-33.33	-32.15		-40.14	-36.33		13.00%
1R69	69	30	-34.2	-31.85	1:07:32	-33.35	-32.20		-40.85	-36.28		12.67%
1CTF	74	42	-38	-35.28	1:37:44	-45.83	-40.94		-51.5	-47.29		15.51%
3MX7	90	44	<i>n/a</i>	<i>n/a</i>		-44.81	-42.32		-56.32	-50.95		20.39%
3NBM	108	56	<i>n/a</i>	<i>n/a</i>		-52.44	-49.51		-53.66	-49.9		0.79%
3MQO	120	68	<i>n/a</i>	<i>n/a</i>		-64.04	-58.84	1:00:00	-62.25	-54.56	1:00:00	<i>no RI</i>
3MRO	142	63	<i>n/a</i>	<i>n/a</i>		-87.38	-82.24		-90.05	-82.32		0.10%
3PNX	160	84	<i>n/a</i>	<i>n/a</i>		-103.04	-96.86		-102.55	-88.06		<i>no RI</i>
3MSE	180	83	<i>n/a</i>	<i>n/a</i>		<i>n/a</i>	<i>n/a</i>		-92.61	-84.60		<i>n/a</i>
3MR7	189	88	<i>n/a</i>	<i>n/a</i>		<i>n/a</i>	<i>n/a</i>		-93.65	-83.93		<i>n/a</i>
3MQZ	215	115	<i>n/a</i>	<i>n/a</i>		<i>n/a</i>	<i>n/a</i>		-104.29	-95.22	2:00:00	<i>n/a</i>
3NO3	238	102	<i>n/a</i>	<i>n/a</i>		<i>n/a</i>	<i>n/a</i>		-122.97	-108.70		<i>n/a</i>
3NO6	248	112	<i>n/a</i>	<i>n/a</i>		<i>n/a</i>	<i>n/a</i>		-133.95	-117.11		<i>n/a</i>
3ON7	280	135	<i>n/a</i>	<i>n/a</i>		<i>n/a</i>	<i>n/a</i>		-116.88	-96.64		<i>n/a</i>

n/a denotes the experimental results are not available.

Table 5: The average energy and average RMSD values achieved from two different variants of GA. The average values are calculated over 50 different runs. The bold-faced values indicate the winner (the lower the better).

Protein details			The state-of-the-art GA [61]				The MH_GA					
			Reported values				Average values					
			MJ model			MJ model		MH model		Gen		
Seq	Size	H	Energy	RMSD	Pop	Gen	Energy	RMSD	Energy	RMSD	Pop	(≤)
2J6A	135	71	-815.82*	16.75	50	20000	-59.72	9.53	-61.40	9.48	50	2500
2HFQ	85	38	-543.17*	12.24	50	20000	-52.13	7.48	-52.72	7.31	50	7000

* the unusual values for MJ energy model.

Shatabda *et al.* [65] used 3D FCC lattice with 20×20 empirical energy matrix by Berrera *et al.* [11]. In fact, we do not have any state-of-the-art results available for similar model to compare free energy level in a straight way. Therefore, we ran the algorithms used in [62] and [65] using the MJ energy model [10] to compare our results. However, the constraint programming based hybrid approach [62] failed to get any solution for most of the large-sized proteins. In such cases, in Table 4, the results are denoted by *n/a*.

In Table 4, we present interaction energy values in two different formats: the global lowest interaction

energy (Column *Best*) and the average (Column *Avg*) of the lowest interaction energies obtained from 50 different runs. In case of the global best energy, our approach outperforms the state-of-the-art approaches in [62, 65] on 9 out of 12 benchmark proteins. However, in case of average energy, our approach outperforms both of the approaches on 10 out of 12 benchmark proteins. Based on the experimental results, the performance hierarchy of the approaches used to validate our MH_GA is shown in Figure 8.



Figure 8: The performance hierarchy among the state-of-the-art approaches and our MH_GA. Our GA outperforms the other to approaches in [62] and [65].

Outcome based on Relative Improvement (RI)

From the Column RI of Table 4, we see that for 2 proteins our GA fail to improve over the state-of-the-art. However, for other 10 proteins it improves the average interaction energy level ranging from 0.10% to 26.58% for different proteins.

Further, in Table 5, we present another two benchmark proteins taken from a GA based approach [61]. From the authors of [61], we tried to get their implemented codes so that we can run that by ourselves. However, we failed to receive any response from the authors. Therefore, we present the reported values. For fair comparison, we compare the results by generation-wise instead of by running-time.

Outcome based on RMSD comparison

We calculate RMSD of a structure that corresponds to the lowest MJ interaction energy for a particular run. The reported RMSD values in Table 6 are the global minimum of 50 runs. In Table 5 and Table 6, the bold-faced RMSD values indicate the winners for the corresponding proteins.

In Table 7, we present corresponding MJ energy values for global minimum RMSD and corresponding RMSD values for global minimum MJ energy values over 50 runs for each proteins on identical settings. The experimental results show that the global minimum energy in our experiment does not produce minimum RMSD value.

5.3. Result Analysis

The MJ energy model actually implicitly bear the characteristic of hydrophobicity. The matrix values present some variations within amino acids of the same class (H or P). A partition algorithm such as 2-

Table 6: The best RMSD values reported, are the best amongst the 50 different runs. The bold-faced values indicate the winner (the lower the better).

Protein details			Local Search [65]		The MH_GA	
<i>Seq</i>	<i>Size</i>	<i>H</i>	<i>MJ guided</i>	<i>HP guided</i>	<i>MJ guided</i>	<i>MH guided</i>
4RXN	54	27	5.74	4.70	4.83	4.76
1ENH	54	19	5.94	4.42	4.75	4.81
4PTI	58	32	6.02	6.18	6.24	6.06
2IGD	61	25	7.38	7.64	6.63	6.53
1YPA	64	38	6.54	5.17	5.52	5.39
1R69	69	30	6.12	4.44	4.76	4.64
1CTF	74	42	6.08	4.72	4.26	4.08
3MX7	90	44	8.17	7.10	7.21	7.20
3NBM	108	56	6.38	5.89	5.64	5.37
3MQO	120	68	6.92	6.44	6.33	6.38
3MRO	142	63	8.76	7.76	7.93	7.64
3PNX	160	84	8.78	7.90	8.04	7.60
3MSE	180	83	<i>n/a</i>	20.24	16.05	16.98
3MR7	189	88	<i>n/a</i>	10.43	9.42	9.36
3MQZ	215	115	<i>n/a</i>	11.21	8.88	9.04
3NO3	238	102	<i>n/a</i>	14.49	11.22	11.70
3NO6	248	112	<i>n/a</i>	13.20	13.88	12.04
3ON7	280	135	<i>n/a</i>	13.19	11.84	11.77

means clustering algorithm easily reveals the H-P partitioning within the MJ model. Given this knowledge, we study the effect of explicitly using hydrophobic property within our GA.

Effect of HP in MH model

Our macro-mutation operator biases the search towards a hydrophobic core by applying a series of diagonal moves and thus achieves improvements in terms of MJ energy values of the output conformations. We implemented three different versions of our genetic algorithm.

1. **MH:** This version is our final algorithm that we described in detail, and used in presenting our main results in Table 4 and in comparing with the state-of-the-art results. To reiterate, this version uses the MJ energy model for search and energy reporting, and hydrophobicity knowledge in the macro-mutation operator that repeatedly applies diagonal moves towards forming a hydrophobic core.
2. **MJ:** This version of our GA uses the MJ energy model for search and energy reporting. This version has macro-mutation operator but not biased by hydrophobic properties of amino acids.
3. **HP:** This version of our GA uses the HP energy model for search. However, we report the energy values of the final conformations returned by the GA in MJ energy model. Note that this version has the hydrophobic core directed macro-mutation operator. This version will show whether HP model is sufficient even when the energy of a conformation is to be in the MJ model.

Table 7: Corresponding MJ energies for global minimum RMSD and corresponding RMSDs for global minimum MJ energies over 50 runs for each proteins

Protein details			Energy corresponds to RMSD						RMSD corresponds to energy					
			HP		MJ		MH		HP		MJ		MH	
<i>Seq</i>	<i>Size</i>	<i>H</i>	<i>rmsd</i>	<i>En</i>	<i>rmsd</i>	<i>En</i>	<i>rmsd</i>	<i>En</i>	<i>En</i>	<i>rmsd</i>	<i>En</i>	<i>rmsd</i>	<i>En</i>	<i>rmsd</i>
4RXN	54	27	4.70	4.24	4.83	-26.68	4.76	-26.02	-12.41	6.30	-37.06	5.91	-36.36	5.99
1ENH	54	19	4.42	-0.67	4.75	-15.21	4.81	-10.8	-10.27	7.26	-38.85	7.68	-38.39	7.14
4PTI	58	32	6.18	-0.36	6.24	-8.03	6.06	-19.16	-6.95	7.00	-32.6	8.09	-35.65	8.62
2IGD	61	25	7.64	4.00	6.63	-18.21	6.53	-19.79	-10.28	9.4	-35.57	9.86	-36.49	8.69
1YPA	64	38	5.17	5.21	5.52	-26.90	5.39	-35.01	-17.1	8.37	-38.45	7.81	-40.14	8.32
1R69	69	30	4.44	3.59	4.76	-21.70	4.64	-22.37	-11.3	5.38	-39.89	7.16	-40.85	6.40
1CTF	74	42	4.72	-3.72	4.26	-32.55	4.08	-44.44	-18.06	7.19	-50.45	5.96	-51.50	5.94
3MX7	90	44	7.10	-0.08	7.21	-42.18	7.20	-50.85	-17.97	8.73	-56.55	10.05	-56.32	9.57
3NBM	108	56	5.89	-5.07	5.64	-35.75	5.37	-36.51	-23.09	8.27	-55.38	6.75	-53.66	7.34
3MQO	120	68	6.44	5.96	6.33	-51.44	6.38	-41.69	-15.47	9.31	-62.65	7.69	-62.25	8.13
3MRO	142	63	7.76	-10.97	7.93	-50.69	7.64	-68.41	-28.63	12.96	-90.56	11.89	-90.05	9.28
3PNX	160	84	7.90	-1.16	8.04	-73.90	7.60	-69.52	-26.79	10.81	-96.98	10.11	-102.55	10.12
3MSE	180	83	20.24	-14.41	16.05	-76.99	16.98	-77.73	-30.4	22.01	-91.02	19.12	-92.61	17.88
3MR7	189	88	10.43	-12.34	9.42	-84.28	9.36	-81.9	-26.99	10.56	-94.93	11.67	-93.65	10.84
3MQZ	215	115	11.21	-5.26	8.88	-98.75	9.04	-92.85	-15.51	11.53	-108.38	10.58	-104.29	10.7
3NO3	238	102	14.49	-14.51	11.22	-112.14	11.7	-100.79	-16.41	14.89	-119.9	13.2	-122.97	13.04
3NO6	248	112	13.2	-8.67	11.88	-120.23	12.06	-116.51	-44.07	13.96	-125.68	14.26	-133.95	13.09
3ON7	280	135	13.19	28.47	11.84	-105.63	11.77	-98.31	-8.59	13.95	-120.16	13.01	-116.88	16.58

From the Column RI in Table 8, we see that MH guided GA improves the average interaction energy level over MJ model, ranging from 0.84% to 5.14% for all benchmark proteins. The improvements are not large in magnitudes but consistently better for all the proteins.

Statistical significance

We know that the lower *p-values* are better. We performed the *t-test* with a confidence interval of 95% (i.e., significance level is 5%) and the results are presented in Table 8. For MJ and MH models, the p-values of all proteins are less than the significance level. However, for HP model, the p-value for 3MSE is below the significance level and for other five sequences those are equal to the significance level. Therefore, the experimental results are statistically significant.

Search progress

To demonstrate the search progress, we periodically find the best energy values obtained so far in each run. For a given period, we then calculate the average energy values obtained for that period over 50 runs. We used a 2 minute time interval. Figure 9 presents the average energy values obtained at each time interval for two different proteins: *4RXN* and *3PNX* are the smallest and largest amongst the 12 benchmark proteins

Table 8: The effect of using HP energy model within a macro-mutation operator. The bold-faced values indicate the winners. The lower the energy value, the better the performance. The t -test was performed with a confidence interval of 95%.

Protein details			Best of 50 runs			Average[p-value] of 50 runs			RI
<i>Seq</i>	<i>Size</i>	<i>H</i>	<i>HP</i>	<i>MJ</i>	<i>MH</i>	<i>HP</i>	<i>MJ(r)</i>	<i>MH(t)</i>	<i>on MJ</i>
4RXN	54	27	-12.41	-37.71	-36.36	-3.54[2.4E-16]	-33.32[5.9E-56]	-33.60 [1.7E-75]	0.84%
1ENH	54	19	-10.27	-37.37	-38.39	-7.29[3.8E-32]	-34.86[1.1E-66]	-35.67 [1.2E-70]	2.32%
4PTI	58	32	-6.95	-35.31	-35.65	-2.81[1.5E-14]	-30.93[3.6E-55]	-31.01 [4.8E-67]	0.26%
2IGD	61	25	-10.28	-36.97	-36.49	-6.75[2.7E-31]	-33.65[3.5E-66]	-33.75 [4.0E-70]	0.30%
1YPA	64	38	-17.1	-39.13	-40.14	-9.90[2.3E-33]	-35.20[6.4E-65]	-36.33 [2.8E-73]	3.21%
1R69	69	30	-11.3	-39.77	-40.85	-4.31[5.6E-19]	-35.43[4.9E-65]	-36.28 [2.5E-68]	2.40%
1CTF	74	42	-18.06	-50.09	-51.5	-10.97[1.1E-32]	-44.98[1.4E-61]	-47.29 [6.8E-70]	5.14%
3MX7	90	44	-17.97	-55.57	-56.32	-11.16[1.9E-31]	-48.46[5.5E-62]	-50.95 [2.6E-70]	5.14%
3NBM	108	56	-23.09	-57.17	-53.66	-15.29[9.8E-36]	-48.47[9.5E-60]	-49.90 [2.6E-70]	2.95%
3MQO	120	68	-15.47	-60.22	-62.25	-6.75[1.7E-18]	-53.00[4.8E-61]	-54.56 [2.4E-66]	2.94%
3MRO	142	63	-28.63	-93.77	-90.05	-18.65[7.2E-31]	-79.32[2.1E-62]	-82.32 [1.6E-67]	3.78%
3PNX	160	84	-26.79	-99.87	-102.55	-18.55[1.2E-34]	-85.64[6.0E-60]	-88.06 [1.3E-60]	2.83%
3MSE	180	83	-30.4	-91.02	-92.61	-13.17[5.0E-21]	-84.47[3.2E-70]	-84.60 [3.6E-69]	0.20%
3MR7	189	88	-26.99	-94.93	-93.65	-5.54[1.4E-06]	-85.70 [4.1E-69]	-83.93[1.9E-36]	<i>non</i>
3MQZ	215	115	-15.51	-108.38	-104.29	6.86[8.7E-08]	-96.58 [1.5E-68]	-95.22[6.7E-64]	<i>non</i>
3NO3	238	102	-16.41	-119.9	-122.97	-2.41[5.1E-02]	-108.68[1.1E-68]	-108.70 [3.3E-65]	0.12%
3NO6	248	112	-44.07	-125.68	-133.95	-12.65[2.0E-11]	-116.31[1.8E-71]	-117.11 [7.0E-67]	0.70%
3ON7	280	135	-8.59	-120.16	-116.88	9.38[7.0E-10]	-104.57 [1.1E-56]	-96.64[2.4E-45]	<i>non</i>

respectively. From both of the charts, we see that the final version of our algorithm MH performs better than the other two versions.

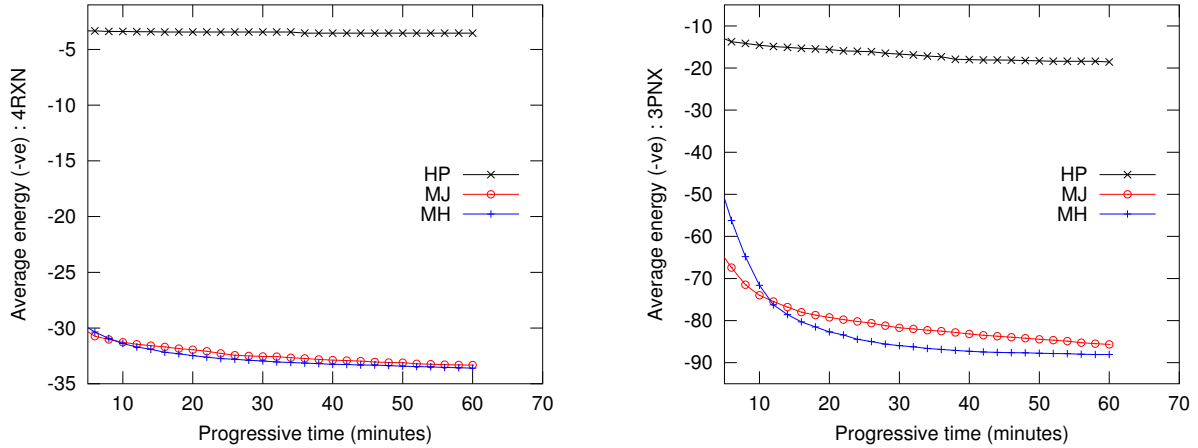


Figure 9: The search progress over a time-span of 60 minutes for proteins *4RXN* and *3PNX* of sequence length 54 and 160 amino acids respectively.

6. Discussion

By encoding the conformation with angular coordinates (ϕ and ψ), our GA might easily be applied in high-resolution PSP. While the minimizing energy function is highly complex (such as molecular dynamics), a simple guidance heuristic—such as hydrophobic property or exposed surface area—could be used to guide the macro-mutation operator. Within GA framework, the macro-mutation operator could be applied optimizing the segments of secondary structures (α -helix and β -sheet).

Our approach can easily divide the whole optimization process into two stages guided by two energy models with different complexities. The macro-mutation operator can be guided by simpler energy models such as distance from hydrophobic core, exposed surface area, hydrophobicity of amino acids, hydrophathy index of the amino acids, and so on. Conversely, the main objective function can be more realistic such as molecular dynamics based energy models. This two-stage optimization will reduce the overall computational complexities. As a result, our framework has a good chance to succeed in more realistic models even for large sized proteins.

7. Conclusion

Our guided macro-mutation in a graded energy based genetic algorithm, ‘MH_GeneticAlgorithm’, is found to be an effective sampling algorithm for the convoluted protein structure space. The strategical switching in between the Miyazawa-Jernigan (MJ) energy and the Hydrophobic-Polar (HP) energy made the proposed algorithm perform better compared to the other state-of-the-art approaches. This is because, while the fine graded MJ energy interaction computation become computationally prohibit, the low resolution HP energy model can effectively sample the search-space towards certain promising directions. In addition, the GA framework was enhanced and made powerful, since it uses not only crossover but also three effective move operators. Further to diversify the population to keep sampling or, exploring the search space effectively, a hydrophobic core-directed macro-mutation operator, twin removal as well as a random-walk algorithm to recover from the stagnation have been applied. To compare the performance of our GA, we have extensively compared with the existing state-of-the-approaches using the benchmark problem available and found our approach to be consistently better and often found significantly better and t-test result in terms of p-values have been provided. For the lattice configuration to be followed, we used 3D face-centered-cube (FCC) lattice model because prediction in the FCC lattice model can yield the densest protein core and the FCC lattice model can provide the maximum degree of freedom as well as the closest resemblance to the real or, high resolution folding within the lattice constraint. This enables the predicted structure to be aligned and hence, migrated to a real protein (prediction) model efficiently for future extensions.

Acknowledgment

Mahmood Rashid and Abdul Sattar would like to express their great appreciation to National ICT Australia. Sumaiya Iqbal and Md Tamjidul Hoque acknowledge the Louisiana Board of Regents through the Board of Regents Support Fund, LEQSF(2013-16)-RD-A-19.

References

References

- [1] H. J. Morowitz, Energy flow in biology, Academic Press, 1968.
- [2] A. Stouthamer, A theoretical study on the amount of ATP required for synthesis of microbial cell material, *Antonie van Leeuwenhoek* 39 (1) (1973) 545–565.
- [3] B. Alberts, A. Johnson, J. Lewis, M. Raff, K. Roberts, P. Walters, The shape and structure of proteins, in: *Mol. Bio. of the Cell*, Fourth Ed, 2002.
- [4] Adam Smith, Protein misfolding, *Nature Reviews Drug Discovery* 426 (6968) (2003) 78–102.
- [5] C. M. Dobson, Protein folding and misfolding, *Nature* 426 (6968) (2003) 884–890.
- [6] Eleanor J. Dodson, Computational biology: protein predictions, *Nature* 450 (7167) (2007) 176–177.
- [7] K. W. Kaufmann, G. H. Lemmon, S. L. DeLuca, J. H. Sheehan, J. Meiler, Practically useful: what the rosetta protein modeling suite can do for you, *Biochemistry* 49 (14) (2010) 2987–2998.
- [8] D. Xu, Y. Zhang, Ab initio protein structure assembly using continuous structure fragments and optimized knowledge-based force field, *Proteins: Structure, Function, and Bioinformatics* 80 (7) (2012) 1715–1735.
- [9] J. Lee, S. Wu, Y. Zhang, Ab initio protein structure prediction, in: *From protein structure to function with bioinformatics*, Springer, 2009, pp. 3–25.
- [10] S. Miyazawa, R. L. Jernigan, Estimation of effective interresidue contact energies from protein crystal structures: quasi-chemical approximation, *Macromolecules* 18 (3) (1985) 534–552.
- [11] M. Berrera, H. Molinari, F. Fogolari, Amino acid empirical contact energy definitions for fold recognition in the space of contact maps, *BMC Bioinformatics* 4 (1) (2003) 8.
- [12] K. F. Lau, K. A. Dill, A lattice statistical mechanics model of the conformational and sequence spaces of proteins, *Macromolecules* 22 (10) (1989) 3986–3997.
- [13] S. Cooper, F. Khatib, A. Treuille, J. Barbero, J. Lee, M. Beenen, A. Leaver-Fay, D. Baker, Z. Popovič, F. Players, Predicting protein structures with a multiplayer online game, *Nature* 466 (7307) (2010) 756–760.
- [14] R. Das, D. Baker, Macromolecular modeling with rosetta, *Annu. Rev. Biochem.* 77 (2008) 363–382.
- [15] R. Wroe, E. Bornberg-Bauer, H. S. Chan, Comparing folding codes in simple heteropolymer models of protein evolutionary landscape: robustness of the superfunnel paradigm, *Biophysical journal* 88 (1) (2005) 118–131.
- [16] M. T. Hoque, M. Chetty, L. S. Dooley, A new guided genetic algorithm for 2d hydrophobic-hydrophilic model to predict protein folding, *IEEE Congress on Evolutionary Computation* 1 (2005) 259–266.

- [17] M. T. Hoque, M. Chetty, A. Lewis, A. Sattar, Twin removal in genetic algorithms for protein structure prediction using low-resolution model, *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 8 (1) (2011) 234–245.
- [18] M. T. Hoque, M. Chetty, A. Lewis, A. Sattar, V. M. Avery, DFS-generated pathways in ga crossover for protein structure prediction, *Neurocomputing* 73 (13) (2010) 2308–2316.
- [19] S. Iqbal, A. Mishra, M. T. Hoque, Improved prediction of accessible surface area results in efficient energy function application, *Journal of Theoretical Biology* 380 (2015) 380–391.
- [20] T. B. Higgs, B. Stantic, M. T. Hoque, A. Sattar, Applying feature-based resampling to protein structure prediction, 2012.
URL http://cs.uno.edu/~tamjid/Papers/2012_FBR_PSP.pdf
- [21] E. Giaquinta, L. Pozzi, An effective exact algorithm and a new upper bound for the number of contacts in the hydrophobic-polar two-dimensional lattice model, *Journal of Computational Biology* 20 (8) (2013) 593–609.
- [22] R. Unger, J. Moult, Genetic algorithms for protein folding simulations, *Journal of molecular biology* 231 (1) (1993) 75–81.
- [23] N. Lesh, M. Mitzenmacher, S. Whitesides, A complete and effective move set for simplified protein folding, in: *Research in Computational Molecular Biology*, ACM, 2003, pp. 188–195.
- [24] M. T. Hoque, M. Chetty, L. S. Dooley, A new guided genetic algorithm for 2D hydrophobic-hydrophilic model to predict protein folding, *IEEE Congress on Evolutionary Computation* 1 (2005) 259–266.
- [25] M. T. Hoque, M. Chetty, A. Sattar, Protein folding prediction in 3D FCC HP lattice model using genetic algorithm, 2007, pp. 4138–4145.
- [26] C. Thachuk, A. Shmygelska, H. H. Hoos, A replica exchange Monte Carlo algorithm for protein folding in the HP model, *BMC bioinformatics* 8 (1) (2007) 342.
- [27] A.-A. Tantar, N. Melab, E.-G. Talbi, A grid-based genetic algorithm combined with an adaptive simulated annealing for protein structure prediction, *Soft Computing-A Fusion of Foundations, Methodologies and Applications* 12 (12) (2008) 1185–1198.
- [28] R. Unger, J. Moult, A genetic algorithm for 3D protein folding simulations, Morgan Kaufmann Publishers, The 5th International Conference on Genetic Algorithms, 1993, p. 581.
- [29] M. T. Hoque, Genetic Algorithm for *ab initio* protein structure prediction based on low resolution models, Ph.D. thesis, Monash University, Australia (Sep. 2007).
- [30] H.-J. Böckenhauer, A. Z. M. D. Ullah, L. Kapsokalivas, K. Steinhöfel, A local move set for protein folding in triangular lattice models, in: *WABI*, Vol. 5251 of *Lecture Notes in Computer Science*, Springer, 2008, pp. 369–381.
- [31] M. T. Hoque, M. Chetty, A. Lewis, A. Sattar, Twin removal in genetic algorithms for protein structure prediction using low-resolution model, *Transactions on Computational Biology and Bioinformatics* 8 (1) (2011) 234–245.
- [32] G. W. Klau, N. Lesh, J. Marks, M. Mitzenmacher, Human-guided tabu search, in: *The Eighteenth National Conference on Artificial Intelligence (AAAI-02)*, 2002.
- [33] C. Blum, Ant colony optimization: Introduction and recent trends, *Physics of Life reviews* 2 (4) (2005) 353–373.

- [34] I. Kondov, R. Berlich, Protein structure prediction using particle swarm optimization and a distributed parallel approach, in: ACM workshop on Biologically inspired algorithms for distributed systems, BADS '11, 2011, pp. 35–42.
- [35] N. Mansour, F. Kanj, H. Khachfe, Particle swarm optimization approach for protein structure prediction in the 3D HP model., *Interdiscip Sci* 4 (3) (2012) 190–200.
- [36] V. Cutello, G. Nicosia, M. Pavone, J. Timmis, An immune algorithm for protein structure prediction on lattice models, *Evolutionary Computation*, IEEE Transactions on 11 (1) (2007) 101–117.
- [37] M. Cebrián, I. Dotú, P. Van Hentenryck, P. Clote, Protein structure prediction on the face centered cubic lattice by local search, in: Proceedings of the 23rd national conference on Artificial intelligence - Volume 1, 2008, pp. 241–246.
- [38] S. Shatabda, M. A. H. Newton, D. N. Pham, A. Sattar, Memory-based local search for simplified protein structure prediction, in: The ACM Conference on Bioinformatics, Computational Biology and Biomedicine (ACM-BCB), ACM, Orlando, FL, USA, 2012.
- [39] B. Maher, A. A. Albrecht, M. Loomes, X.-S. Yang, K. Steinhöfel, A firefly-inspired method for protein structure prediction in lattice models, *Biomolecules* 4 (1) (2014) 56–75.
- [40] M. Mann, S. Will, R. Backofen, CPSP-tools – Exact and complete algorithms for high-throughput 3D lattice protein studies, *Bmc Bioinformatics* 9 (1) (2008) 230.
- [41] I. Dotú, M. Cebrián, P. Van Hentenryck, P. Clote, On lattice protein structure prediction revisited, *IEEE Trans. on Comp. Biology and Bioinformatics* 8 (6) (2011) 1620 – 32.
- [42] N. Krasnogor, B. P. Blackburne, E. K. Burke, J. Hirst, Multimeme algorithms for protein structure prediction, *Parallel Problem Solving from Nature PPSN VII* 2439 (2002) 769 – 778.
- [43] D. A. Pelta, N. Krasnogor, Multimeme algorithms using fuzzy logic based memes for protein structure prediction, in: Recent advances in memetic algorithms, Springer, 2005, pp. 49–64.
- [44] M. K. Islam, Memetic approach for prediction of low resolution protein structures using lattice models, Ph.D. thesis, Monash University, Victoria, Australia (2011).
- [45] M. K. Islam, M. Chetty, Novel memetic algorithm for protein structure prediction, in: Australasian Conference on Artificial Intelligence, 2009, pp. 412–421.
- [46] M. K. Islam, M. Chetty, A. Z. M. D. Ullah, K. Steinhöfel, A memetic approach to protein structure prediction in triangular lattices, in: *ICONIP* (1), 2011, pp. 625–635.
- [47] M. K. Islam, M. Chetty, M. Murshed, Conflict resolution based global search operators for long protein structures prediction, in: *ICONIP* (1), 2011, pp. 636–645.
- [48] M. K. Islam, M. Chetty, Clustered memetic algorithm for protein structure prediction, in: IEEE Congress on Evolutionary Computation, 2010, pp. 1–8.
- [49] M. K. Islam, M. Chetty, M. Murshed, Novel local improvement techniques in clustered memetic algorithm for protein structure prediction, in: IEEE Congress on Evolutionary Computation, 2011, pp. 1003–1011.
- [50] M. K. Islam, M. Chetty, Clustered memetic algorithm with local heuristics for ab initio protein structure prediction, *Evolutionary Computation*, IEEE Transactions on 17 (4) (2013) 558–576.

- [51] M. A. Rashid, M. T. Hoque, M. A. H. Newton, D. Pham, A. Sattar, A new genetic algorithm for simplified protein structure prediction, in: *Advances in Artificial Intelligence*, Springer Berlin Heidelberg, 2012, pp. 107–119.
- [52] S. Shatabda, M. A. H. Newton, M. A. Rashid, A. Sattar, An efficient encoding for simplified protein structure prediction using genetic algorithms, in: *Evolutionary Computation (CEC), 2013 IEEE Congress on*, 2013, pp. 1217–1224.
- [53] M. A. Rashid, M. A. H. Newton, M. T. Hoque, S. Shatabda, D. Pham, A. Sattar, Spiral search: a hydrophobic-core directed local search for simplified PSP on 3D FCC lattice, *BMC Bioinformatics* 14 (Suppl 2) (2013) S16.
- [54] M. A. Rashid, M. A. H. Newton, M. T. Hoque, A. Sattar, A local search embedded genetic algorithm for simplified protein structure prediction., in: *IEEE Congress on Evolutionary Computation*, IEEE, 2013, pp. 1091–1098.
- [55] M. A. Rashid, M. T. Hoque, M. A. H. Newton, A. Sattar, Collaborative parallel local search for simplified protein structure prediction., in: *12th IEEE International Symposium on Parallel and Distributed Processing with Applications, ISPA*, IEEE Computer Society, 2013, pp. 966–973.
- [56] C. Kern, L. Liao, Lattice models with asymmetric propensity matrices for locationally informed protein structure prediction, *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)* (2013) 90–93.
- [57] A. Dal Palù, A. Dovier, F. Fogolari, Constraint logic programming approach to protein structure prediction, *BMC bioinformatics* 5 (1) (2004) 186.
- [58] A. Dal Palù, A. Dovier, E. Pontelli, Heuristics, optimizations, and parallelism for protein structure prediction in clp (fd), *Proceedings of the 7th international conference on Principles and practice of declarative programming* (2005) 230–241.
- [59] A. Dal Palù, A. Dovier, F. Fogolari, E. Pontelli, Exploring protein fragment assembly using clp, *Proceedings of the international joint conference on Artificial Intelligence* 3 (2011) 2590–2595.
- [60] L. Kapsokalivas, X. Gan, A. A. Albrecht, K. Steinhöfel, Population-based local search for protein folding simulation in the MJ energy model and cubic lattices, *Comput. Biol. Chem.* 33 (4) (2009) 283–294.
- [61] S. R. D. Torres, D. C. B. Romero, L. F. N. Vasquez, Y. J. P. Ardila, A novel *ab-initio* genetic-based approach for protein folding prediction, in: *Proceedings of the 9th annual conference on Genetic and evolutionary computation, GECCO '07*, ACM, 2007, pp. 393–400.
- [62] A. D. Ullah, K. Steinhöfel, A hybrid approach to protein folding problem integrating constraint programming with local search, *BMC bioinformatics* 11 (Suppl 1) (2010) S39.
- [63] A. Dal Palù, E. Pontelli, A. Dovier, A constraint solver for discrete lattices, its parallelization, and application to protein structure prediction, *Software-Practice and Experience* 37 (13) (2007) 1405.
- [64] A. D. Ullah, L. Kapsokalivas, M. Mann, K. Steinhöfel, Protein folding simulation by two-stage optimization, in: *Computational Intelligence and Intelligent Systems*, Vol. 51, Springer Berlin Heidelberg, Berlin, Heidelberg, 2009, pp. 138–145.
- [65] S. Shatabda, M. A. H. Newton, A. Sattar, Mixed heuristic local search for protein structure prediction, in: *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.
- [66] M. A. Rashid, M. A. H. Newton, M. T. Hoque, A. Sattar, Mixing energy models in genetic algorithms for on-lattice protein structure prediction, *BioMed Research International* (2013) 15.
- [67] C. B. Anfinsen, The principles that govern the folding of protein chains, *Science* 181 (4096) (1973) 223–230.

- [68] C. Levinthal, Are there pathways for protein folding?, *Journal of Medical Physics* 65 (1) (1968) 44–45.
- [69] E. Alm, D. Baker, Matching theory and experiment in protein folding, *Current opinion in structural biology* 9 (2) (1999) 189–196.
- [70] D. Baker, A surprising simplicity to protein folding, *Nature* 405 (6782) (2000) 39–42.
- [71] S. Istrail, F. Lam, et al., Combinatorial algorithms for protein folding in lattice models: a survey of mathematical results, *Communications in Information & Systems* 9 (4) (2009) 303–346.
- [72] J. M. Bahi, C. Guyeux, K. Mazouzi, L. Philippe, Computational investigations of folded self-avoiding walks related to protein folding, *Computational biology and chemistry* 47 (2013) 246–256.
- [73] T. C. Hales, A proof of the Kepler conjecture, *The Annals of Mathematics* 162 (3) (2005) 1065–1185.
- [74] J. H. Holland, *Adaptation in natural and artificial system: an introduction with application to biology, control and artificial intelligence*, Ann Arbor, University of Michigan Press.
- [75] M. T. Hoque, M. Chetty, L. S. Dooley, Non-isomorphic coding in lattice model and its impact for protein folding prediction using genetic algorithm, in: *IEEE Symposium on Computational Intelligence and Bioinformatics and Computational Biology*, IEEE, 2006, pp. 1–8.
- [76] K. Yue, K. A. Dill, Sequence-structure relationships in proteins and copolymers, *Physical Review E* 48 (3) (1993) 2267.
- [77] T. Higgs, B. Stantic, M. T. Hoque, A. Sattar, Refining genetic algorithm twin removal for high-resolution protein structure prediction, in: *IEEE Congress on Evolutionary Computation (CEC)*, IEEE, 2012, pp. 1–8.
- [78] M. T. Hoque, M. Chetty, L. S. Dooley, Generalized schemata theorem incorporating twin removal for protein structure prediction, in: *Pattern Recognition in Bioinformatics*, Springer, 2007, pp. 84–97.
- [79] M. A. Rashid, S. Shatabda, M. A. H. Newton, M. T. Hoque, D. N. Pham, A. Sattar, Random-walk: a stagnation recovery technique for simplified protein structure prediction., in: *BCB, ACM*, 2012, pp. 620–622.